

# Markov Decision Processes

# Last Time

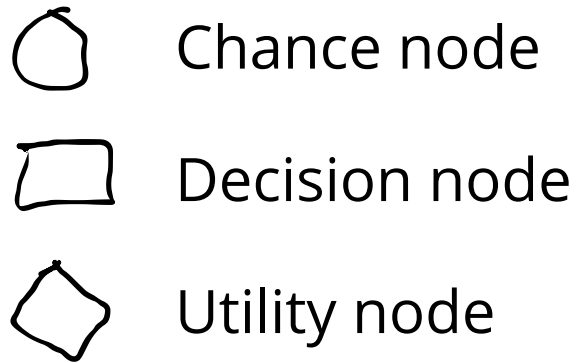
- What does "Markov" mean in "Markov Process"?

# Guiding Questions

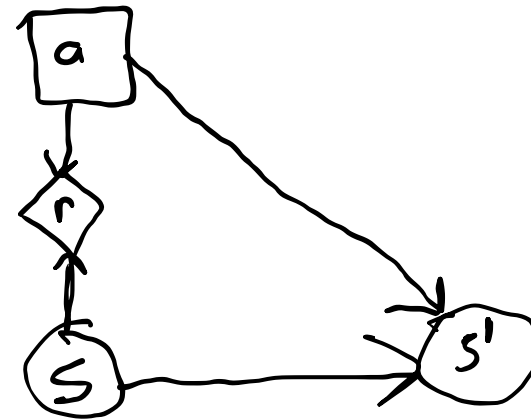
- What is a **Markov decision process**?
- What is a **policy**?
- How do we **evaluate** policies?

# Decision Networks and MDPs

## Decision Network



## MDP Dynamic Decision Network



## MDP Optimization problem

$$\text{maximize } \mathbb{E} \left[ \sum_{t=0}^{\infty} r_t \right]$$

Not well formulated!  
Infinite

# Finite MDP Objectives

1. Finite time

$$\mathbb{E} \left[ \sum_{t=0}^T r_t \right]$$

2. Average reward

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[ \frac{1}{n} \sum_{t=0}^n r_t \right]$$

3. Discounting

$$\mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r_t \right]$$

discount  $\gamma \in [0, 1)$

typically 0.9, 0.95, 0.99

if  $\underline{r} \leq r_t \leq \bar{r}$

4. Terminal States

Infinite time, but a terminal state (no reward, no leaving) is always reached with probability 1.

then

$$\frac{\underline{r}}{1 - \gamma} \leq \sum_{t=0}^{\infty} \gamma^t r_t \leq \frac{\bar{r}}{1 - \gamma}$$

# MDP "Tuple Definition"

$(S, A, T, R, \gamma)$  (and  $b$  and/or  $S_T$  in some contexts)

- $S$  (state space) - set of all possible states
 

$\{1, 2, 3\}$ 
 $(x, y) \in \mathbb{R}^2$ 
 $\{0, 1\} \times \mathbb{R}^4$
- $A$  (action space) - set of all possible actions
 

$\{\text{healthy, pre-cancer, cancer}\}$ 
 $(s, i, r) \in \mathbb{N}^3$
- $T$  (transition distribution) - explicit or implicit ("generative") model of how the state changes
 

$\{1, 2, 3\}$ 
 $\mathbb{R}^2$ 
 $\{0, 1\} \times \mathbb{R}^2$ 
 $\{\text{test, wait, treat}\}$ 
 $T(s' \mid s, a)$
- $R$  (reward function) - maps each state and action to a reward
 

$R(s, a)$  or  
 $R(s, a, s')$
- $\gamma$ : discount factor
 

$s', r = G(s, a)$
- $b$ : initial state distribution
- $S_t$ : set of terminal states

# MDP Example

Imagine it's a cold day and you're ready to go to work. You have to decide whether to bike or drive.

- If you drive, you will have to pay \$15 for parking; biking is free.
- On 1% of cold days, the ground is covered in ice and you will crash if you bike, but you can't discover this until you start riding. After your crash, you limp home with pain equivalent to losing \$100.

# Policies and Simulation

- A *policy*, denoted with  $\pi(a_t \mid s_t)$ , is a conditional distribution of actions given states.
- $a_t = \pi(s_t)$  is used as shorthand when a policy is deterministic.
- When a policy is combined with an MDP, it becomes a Markov stochastic process with

$$P(s' \mid s) = \sum_{a_t} T(s' \mid s, a_t) \pi(a_t \mid s_t)$$



# Break

- Suggest a policy that you think is optimal for the icy day problem

# Policy Evaluation

Naive Policy Evaluation not on Exam

# Monte Carlo Policy Evaluation

- Running a large number of simulations and averaging the accumulated reward is called *Monte Carlo Evaluation*

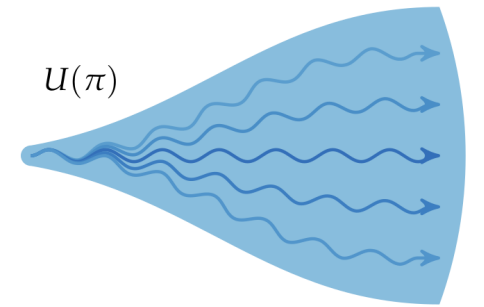
Let  $\tau = (s_0, a_0, r_0, s_1, \dots, s_T)$  be a *trajectory* of the MDP

$$U(\pi) \approx \frac{1}{m} \sum_{i=1}^m R(\tau^{(i)})$$

also an R.V.

$$U(\pi) \approx \bar{u}_m = \frac{1}{m} \sum_{i=1}^m \hat{u}^{(i)}$$

where  $\hat{u}^{(i)}$  is generated by a rollout simulation



How can we quantify the accuracy of  $\bar{u}_m$ ?

# Value Function-Based Policy Evaluation

# Guiding Questions

- What is a **Markov decision process**?
- What is a **policy**?
- How do we **evaluate** policies?